

Exact solutions for the selection–mutation equilibrium in the Crow–Kimura evolutionary model

Yuri S. Semenov¹, Artem S. Novozhilov^{2,*}

¹*Applied Mathematics–1, Moscow State University of Railway Engineering,
Moscow 127994, Russia*

²*Department of Mathematics, North Dakota State University, Fargo, ND 58108, USA*

Abstract

We reformulate the eigenvalue problem for the selection–mutation equilibrium distribution in the case of a haploid asexually reproduced population in the form of an equation for an unknown probability generating function of this distribution. The special form of this equation in the infinite sequence limit allows us to obtain analytically the steady state distributions for a number of particular cases of the fitness landscape. The general approach is illustrated by examples; theoretical findings are compared with numerical calculations.

Keywords: Selection–mutation equilibrium, quasispecies model, Crow–Kimura model, error threshold, single peaked landscape

AMS Subject Classification: Primary: 92D15; 92D25; Secondary: 15A18

1 Introduction

Selection and mutation are two main evolutionary forces that shape (together with recombination and genetic drift) the life histories of evolving populations. The mathematical theory of selection–mutation models is deep and elaborate (e.g., [7]), covering various biological assumptions, such as different mutation schemes, consequence of ploidy, mating systems, heterogeneous environment, etc. The scope of the theory notwithstanding, even the simplest possible formulation of a multi locus mutation–selection model in the case of an asexually reproduced haploid population still presents mathematical challenges. The goal of our manuscript is to introduce an approach that allows, at least in some special cases, a relatively straightforward derivation of the steady state distribution for such model.

Consider a haploid one locus asexually reproduced population with $N + 1$ alleles. Let $p_i = p_i(t)$ be the frequency of the i -th allele at time t , $i = 0, \dots, N$; the corresponding Malthusian fitness is denoted m_i . Let μ_{ij} denote the mutation rate of allele j to allele i . Then, assuming

*Corresponding author: artem.novozhilov@ndsu.edu

that the reproduction events and mutations are separated, we end up with a nonlinear system of ordinary differential equations (e.g., [1, 8]) of the form

$$\dot{p}_i = (m_i - \bar{m})p_i + \sum_{j=0}^N \mu_{ij}p_j, \quad i = 0, \dots, N, \quad (1.1)$$

where $\mu_{ii} = -\sum_{j \neq i} \mu_{ij}$, and $\bar{m} = \sum_{j=0}^N m_j p_j$ is the mean population fitness. In the matrix form system (1.1) reads

$$\dot{\mathbf{p}} = (\mathbf{M} - \bar{m}\mathbf{I})\mathbf{p} + \mathbf{M}\mathbf{p}, \quad (1.2)$$

where $\mathbf{M} = \text{diag}(m_0, \dots, m_N)$ is the fitness landscape, $\mathbf{M} = (\mu_{ij})$ is the mutation matrix, $\mathbf{p} = (p_0, \dots, p_N)^\top$, and \mathbf{I} is the identity matrix. Model (1.2) is invariant with respect to rescalings of the Malthusian fitnesses as $\tilde{m}_i = m_i + \tilde{m}$ for any constant \tilde{m} and does not allow for lethal mutations.

A basic fact about model (1.2) is as follows. Given that the matrix \mathbf{M} is irreducible, then there exists a unique globally stable positive selection–mutation equilibrium $\lim_{t \rightarrow \infty} \mathbf{p}(t) = \hat{\mathbf{p}}$ that can be found as the normalized eigenvector of the matrix $\mathbf{M} + \mathbf{M}$ corresponding to the strictly dominant eigenvalue $\lambda = \hat{\bar{m}} = \sum_{j=0}^N m_j \hat{p}_j$ (e.g., [7, 21]). In Eigen’s theory of the origin of life, which is equivalent to the selection–mutation approach in haploid populations [24], the eigenvector $\hat{\mathbf{p}}$, which describes the equilibrium frequencies of the self-replicating macromolecules, was called the quasispecies [9, 10, 11].

To obtain further theoretical results on the form of the dominant eigenvalue λ and/or selection–mutation equilibrium $\hat{\mathbf{p}}$, additional assumptions on the form of the fitness landscape \mathbf{M} and/or mutation matrix \mathbf{M} are necessary. For example, a profitable way is to neglect the reverse mutations, i.e., put (e.g., [23])

$$\begin{aligned} \mu_{ij} &> 0, \quad i > j, \\ \mu_{ij} &= 0, \quad i < j, \end{aligned}$$

or even more restrictive (e.g., [22])

$$\begin{aligned} \mu_{j+1,j} &> 0, \\ \mu_{ij} &= 0, \quad i \neq j, i \neq j+1. \end{aligned}$$

A less restrictive assumption is to assume that allele j can mutate only to the neighbors $j-1$ and $j+1$, other mutations are prohibited [1, 3, 13]. In the last case, assuming additionally that

$$\begin{aligned} \mu_{j-1,j} &= \mu j, \\ \mu_{j+1,j} &= \mu(N-j), \\ \mu_{jj} &= -\mu N, \end{aligned} \quad (1.3)$$

for some constant $\mu > 0$, it is possible to write down an explicit solution for the equilibrium frequencies \hat{p}_i in the case of an additive or Fujiyama fitness landscape, defined as $m_j = -Mj$ for some constant $M > 0$ [3, 14, 16] (we reproduce this solution, using our method, below in Example 3.6).

We stress that no simple analytical expressions for the components of $\hat{\mathbf{p}}$ are known for a general fitness landscape \mathbf{M} and the mutation scheme (1.3); even in the case of a single or

sharply peaked landscape, defined as $\mathbf{M} = \text{diag}(m_0, 0, \dots, 0)$, $m_0 > 0$, there exists no analytical solution.

The mutation scheme (1.3) naturally arises if one changes the point of view from a one locus $N + 1$ allele population to a biallelic N locus haploid population, which is scrutinized in the quasispecies theory [1, 6, 10, 15]. To this end, consider a population of sequences, each sequence consists of N sites (loci), and each site can be either in 0 or 1 state (two alleles per each locus). Let $\mu > 0$ denote the mutation rate per site per sequence per time unit, such that 0s mutates to 1s and 1s mutates to 0s with the same rate $\mu > 0$. Also assume that fitness is determined by the number of 1s in the sequence such that we have $N + 1$ different sequence classes with fitnesses m_0, m_1, \dots, m_N (this is sometimes called permutation invariant or symmetric fitness landscape). Then the dynamics of the frequencies of classes are determined by model (1.2) with (1.3), which is often called in the literature the paramuse or Crow–Kimura quasispecies model [1].

It was shown that model (1.2), (1.3) is equivalent to a so-called Ising quantum chain (e.g., [3]). This fact allowed obtaining a number of analytical results about the mean population fitness $\lambda = \hat{m}$ (and for some other population averages) in the selection–mutation equilibrium when the sequence length approaches infinity ($N \rightarrow \infty$) under an appropriate scaling of the model parameters [3].

A similar infinite sequence point of view was taken in [13] (see also [2] for a more recent generalization), where a maximum principle for the mean population fitness was formulated. For our needs the maximum principle can be formulated as follows (we note that a more general case is treated in [13]). Assume that $m_i = Nr_i = Nr(x_i)$, $x_i = i/N \in [0, 1]$ and define $g(x) = \mu(1 - 2\sqrt{x(1-x)})$. Then the scaled equilibrium fitness $\hat{r} = \hat{m}/N$ is given by

$$\hat{r} \approx \hat{r}_\infty = \sup_{x \in [0,1]} (r(x) - g(x)). \quad (1.4)$$

Using the approximate expression for the mean fitness in (1.4) it is possible to obtain expressions for other averages such as the variance per site of fitness and of distance from the fittest class [13]. However, no attempt was made in [13] to obtain analytical expressions for $\hat{\mathbf{p}}_\infty$ (we use index ∞ throughout the text to denote the expressions in the infinite sequence limit).

An exact integral representation of $\hat{\mathbf{p}}_\infty$ for the infinite sequence length for the mutation scheme (1.3) and the single peaked landscape was written in [12], however, transparent analytical expressions were obtained only for \hat{r}_∞ and $\hat{p}_{\infty,0}$. In [18] a full solution for $\hat{p}_{\infty,i}$ was written down by disregarding the reverse mutations from class j to class $j-1$. The same solution was rigorously obtained in [6] together with estimates of the speed of convergence. In [17], using the Hamilton–Jacobi formalism, a general solution for $\hat{\mathbf{p}}_\infty$ depending on an arbitrary scaled fitness landscape $r(x)$ is suggested in an integral form. This general solution is, however, not straightforward to apply to obtain relatively simple analytical expressions for the selection–mutation equilibrium frequencies $\hat{p}_{\infty,i}$. Therefore, we conclude that there exists no simple general way to find the quasispecies distribution $\hat{\mathbf{p}}_\infty$ even under the simplifying assumption of the infinite sequence length.

The goal of the present text is to suggest a straightforward way of calculating the selection–mutation equilibrium for the model (1.2), (1.3) in the infinite sequence length limit that leads to transparent analytical expressions at least for some particular fitness landscapes (for an extensive background for the current work we refer to [6] and [19]).

2 A general approach to solve for the selection–mutation equilibrium

Our goal is to find approximations for the dominant eigenvalue $\lambda = \hat{m}$ and the corresponding normalized eigenvector $\hat{\mathbf{p}}$ (such that $\sum_{i=0}^N \hat{p}_i = 1$ and $\hat{p}_i \geq 0$ for any i) of the eigenvalue problem

$$(\mathbf{M} + \mu \mathbf{Q})\hat{\mathbf{p}} = \lambda \hat{\mathbf{p}}, \quad (2.1)$$

where $\mathbf{M} = \text{diag}(m_0, \dots, m_N)$, $\mu > 0$, and

$$\mathbf{Q} = \begin{bmatrix} -N & 1 & 0 & 0 & \dots & \dots & 0 \\ N & -N & 2 & 0 & \dots & \dots & 0 \\ 0 & N-1 & -N & 3 & \dots & \dots & 0 \\ 0 & 0 & N-2 & -N & \dots & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \dots & 2 & -N & N \\ 0 & 0 & \dots & \dots & 0 & 1 & -N \end{bmatrix}.$$

We note that at the equilibrium both \hat{m} and $\hat{\mathbf{p}}$ are the functions of the mutation rate μ : $\hat{m} = \hat{m}(\mu)$, $\hat{\mathbf{p}} = \hat{\mathbf{p}}(\mu)$. Together with the matrix \mathbf{Q} we also consider a linear differential operator

$$\mathcal{Q}: P(s) \longrightarrow (1 - s^2)P'(s) - N(1 - s)P(s), \quad (2.2)$$

acting on the $(N + 1)$ -dimensional vector space of polynomials of degree less or equal N . By direct calculations, matrix \mathbf{Q} is the matrix of \mathcal{Q} in the standard basis $\{1, s, \dots, s^N\}$. For any $\mathbf{m} = (m_0, \dots, m_N) \in \mathbf{R}^{N+1}$ and $P(s) = \sum_{i=0}^N p_i s^i$ we introduce the notation

$$\mathbf{m} \circ P(s) = \sum_{i=0}^N m_i p_i s^i.$$

Then problem (2.1) can be rewritten for the unknown probability generating function $P(s)$ as

$$\mathbf{m} \circ P(s) + \mu \mathcal{Q}P(s) = \hat{m} P(s), \quad (2.3)$$

where $\hat{m} = \mathbf{m} \circ P(1)$, or

$$\mathbf{m} \circ P(s) + \mu(1 - s^2)P'(s) - \mu N(1 - s)P(s) = \hat{m} P(s), \quad (2.4)$$

with the normalization condition $P(1) = 1$. Inasmuch as problem (2.4) is equivalent to (2.1) then, due to the Perron–Frobenius theorem, there exists a unique solution $P(s)$ satisfying $P(1) = 1$. There is little hope to be able to solve equation (2.4) explicitly (one such example, well known in the literature, is given below, see Example 3.6). It is possible, however, to find approximations of the quantities of interest in the case $N \rightarrow \infty$ at least for some particular fitness landscapes under some additional assumptions.

To formulate the general approach, we introduce the notations

$$\mathbf{r} = \frac{\mathbf{m}}{N}, \quad \hat{r} = \frac{\hat{m}}{N}.$$

After dividing by N equation (2.4) takes the form

$$\mathbf{r} \circ P(s) + \frac{\mu}{N}(1-s^2)P'(s) - \mu(1-s)P(s) = \hat{\tau} P(s).$$

Now we make the following assumptions:

$\mathcal{H}1$: There exists the limit

$$\lim_{N \rightarrow \infty} P(s) = P_\infty(s).$$

The distribution $\hat{p}_{\infty,i} = \hat{p}_i$, $i = 0, 1, 2, \dots$, corresponding to $P_\infty(s) = \sum_{i=0}^{\infty} \hat{p}_i s^i$, will be called the limit distribution.

$\mathcal{H}2$:

$$\lim_{N \rightarrow \infty} \frac{\mu}{N}(1-s^2)P'(s) = 0.$$

$\mathcal{H}3$: For some limit operator \mathbf{r}_∞

$$\lim_{N \rightarrow \infty} \mathbf{r} \circ P(s) = \mathbf{r}_\infty \circ P_\infty(s).$$

Remark 2.1. The assumption $\mathcal{H}2$ is a formal consequence of $\mathcal{H}1$ and $\mathcal{H}3$, but we decided to keep it in the list because it gives the main idea of the suggested method.

If $\mathcal{H}1$ – $\mathcal{H}3$ hold then problem (2.4) is reduced to a nonlinear functional equation with respect to the unknown probability generating function $P_\infty(s)$,

$$-\mu(1-s)P_\infty(s) + \mathbf{r}_\infty \circ P_\infty(s) = \hat{\tau}_\infty P_\infty(s), \quad (2.5)$$

with the conditions

$$\mathbf{r}_\infty \circ P(1) = \hat{\tau}_\infty, \quad P_\infty(1) = 1. \quad (2.6)$$

Problem (2.5)–(2.6) can be effectively solved at least for some simple choices of the fitness landscape \mathbf{M} (see the next section for representative examples). A careful scrutiny of validity of the obtained solutions requires a deeper analysis of the convergence of the eigenvalue \hat{m} and the quasispecies $\hat{\mathbf{p}}$ when $N \rightarrow \infty$. This, for instance, can be done with the help of parametric solutions to (2.1) introduced in [6] (in Appendix A.1 we outline the approach used in [6], in Appendix A.2 we provide an alternative parametric solution approach, which is illustrated by applying it to Example 3.8). Notwithstanding these concerns, a formal solution of (2.5)–(2.6) is of significant value, because, as the examples show, the found solutions closely approximate, even for moderate values of N , numerical solutions of (2.1).

In a sense our approach is a generalization of the so-called random variable technique (e.g., [4]), and assumption $\mathcal{H}2$ implies that we disregard all the mutations from j to $j-1$ classes, similarly to [22, 23]. To see this, consider the mutation scheme of the form

$$\begin{aligned} \mu_{j-1,j} &= \mu_1 j, \\ \mu_{j+1,j} &= \mu_2 (N-j), \\ \mu_{jj} &= -\mu_1 j - \mu_2 (N-j), \end{aligned} \quad (2.7)$$

where $\mu_1 \neq \mu_2$. Then, as it can be directly checked, operator \mathcal{Q} takes the form

$$\mathcal{Q}: P(s) \longrightarrow \mu_2(s-1) \left(N - s \frac{d}{ds} \right) P(s) + \mu_1(1-s)P'(s).$$

After dividing by N and formally taking the limit, the only term that is left is

$$-\mu_2(1-s)P(s),$$

which agrees with (2.5) and shows that for the limit equation the rate of backward mutations μ_1 is neglected.

To conclude this section, we suggest the following approach: Solve problem (2.5)–(2.6) for each μ . For the found solution $P_\infty(s)$ check $\mathcal{H}1$ – $\mathcal{H}3$. If the hypotheses do not hold then notice that $P_\infty(s) \equiv 0$ solves (2.5)–(2.6), with $\hat{r}_\infty(\mu) \equiv 0$. This corresponds to the delocalization phenomenon of the quasispecies distribution, or, in terms of the Ising model, the phase transition, which was called the error threshold in the quasispecies theory. A number of examples illustrating this approach are given in the following section.

3 Examples of the steady state distributions

Here we present several examples of the selection–mutation equilibrium for known and new fitness landscapes; we also compare the analytical results with the numerical calculations. Two well known cases are Example 3.1, treated in [6, 12, 18], and Example 3.6, treated originally in [3, 14, 16]. We present full solutions in these two cases to demonstrate how our approach works. Example 3.8 was discussed and partially analyzed in [6], however, no derivation for the steady state distribution \hat{p}_∞ was provided; here we present all the details. Other examples are new and are not treated anywhere else to the best of our knowledge.

Example 3.1 (Single peaked landscape). We start with a testbed (both numerical [20] and analytical [12]) for the quasispecies model, which was called the single or sharply peaked landscape.

Let

$$\mathbf{r}_\infty = (1, 0, \dots, 0, \dots).$$

Then

$$\mathbf{r}_\infty \circ P_\infty(s) = \sum_{i=0}^{\infty} r_i \hat{p}_i t^i = \hat{p}_0 = P_\infty(0).$$

Equation (2.5) takes the form

$$-\mu(1-s)P_\infty(s) + P_\infty(0) = \hat{r}_\infty P_\infty(s).$$

Plug $s = 0$ in the last expression and find

$$-\mu P_\infty(0) + P_\infty(0) = \hat{r}_\infty P_\infty(0).$$

Assuming that $P_\infty(0) \neq 0$ we find

$$\hat{r}_\infty(\mu) = 1 - \mu,$$

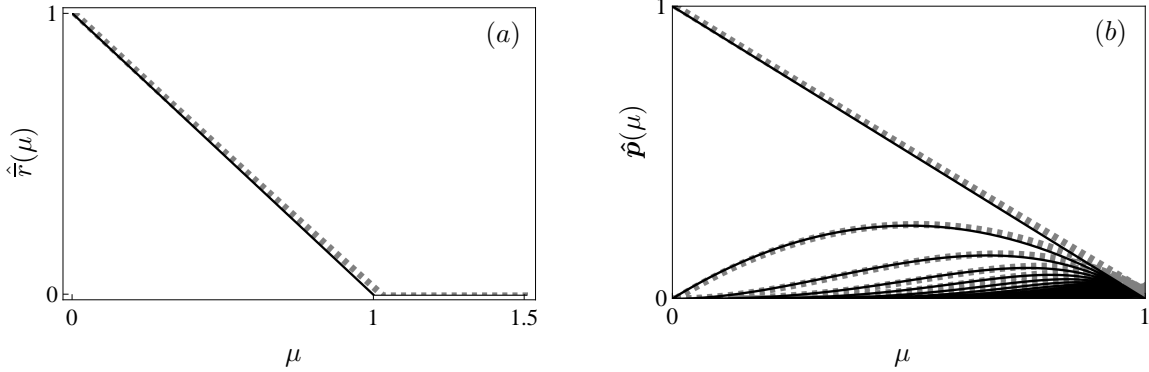


Figure 1: Comparison of numerical calculations for the single peaked landscape $\mathbf{r} = (1, 0, \dots, 0)$ with $N = 50$ with the theoretical predictions of Example 3.1. The black solid lines are the exact solutions for the case $N \rightarrow \infty$ and the grey dashed lines are numerical computations. (a) The mean population fitness versus the mutation rate. (b) The selection–mutation equilibrium $\hat{\mathbf{p}}$ versus the mutation rate. After $\mu \geq 1$ the quasispecies distribution becomes degenerate (binomial)

and hence, using the condition $P_\infty(1) = 1$,

$$P_\infty(s) = \frac{1 - \mu}{1 - \mu s},$$

which is the probability generating function of the geometric distribution with the parameter μ . Therefore the limit distribution is geometric

$$\hat{p}_{\infty,i} = (1 - \mu)\mu^i, \quad i = 0, 1, \dots$$

We see, in view of $\mathcal{H}1$ – $\mathcal{H}3$, that the discussion above holds only for $\mu < 1$, therefore at $\mu = 1$ the structure of the limit distribution abruptly changes and we obtain the solution $P_\infty(s) = 0$ for $\mu \geq 1$. The mean population fitness has the form shown in Fig. 1a. This abrupt change in the quasispecies distribution was called by Eigen et al. the error threshold [5] (see also [6, 13] for an extensive discussion of this notion).

Example 3.2. As a second example, consider a slight generalization of the single peaked landscape in the form

$$\mathbf{r}_\infty = (2, 1, 0, \dots, 0, \dots).$$

Then $\mathbf{r}_\infty \circ P_\infty(s) = 2P_\infty(0) + P'_\infty(0)s$, and (2.5) takes the form

$$-\mu(1 - s)P_\infty(s) + 2P_\infty(0) + P'_\infty(0)s = \hat{r}_\infty P(s).$$

Plugging $s = 0$ yields

$$-\mu P_\infty(0) + 2P_\infty(0) = \hat{r}_\infty P_\infty(0).$$

Assuming $P_\infty(0) \neq 0$ we find

$$\hat{r}_\infty = 2 - \mu,$$

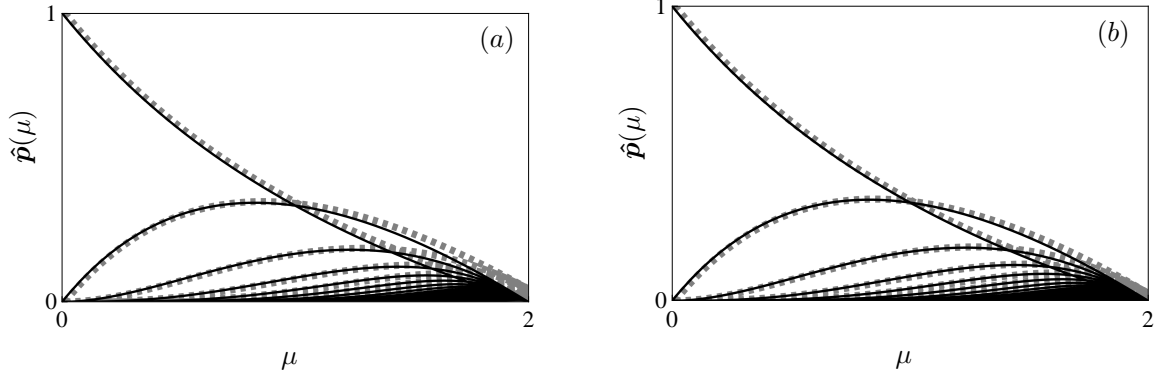


Figure 2: Comparison of numerical calculations for the fitness landscape in Example 3.2 with the theoretical predictions. The black solid lines are the exact solutions for the case $N \rightarrow \infty$ and the grey dashed lines are numerical computations. (a) $N = 50$; (b) $N = 100$

and hence

$$-\mu(1-s)P_\infty(s) + 2P_\infty(0) + P'_\infty(0)s = (2-\mu)P_\infty(s),$$

or

$$\mu s P_\infty(s) + 2P_\infty(0) + P'_\infty(0) = 2P_\infty(s). \quad (3.1)$$

After differentiating the last expression with respect to s and plugging $s = 0$ we find $\mu P_\infty(0) = P'_\infty(0)$. Plugging this into (3.1) implies

$$P_\infty(s) = P_\infty(0) \frac{2 + \mu s}{2 - \mu s},$$

and finally, using the condition $P_\infty(1) = 1$, we obtain

$$P_\infty(s) = \frac{(2-\mu)(2+\mu s)}{(2+\mu)(2-\mu s)} = \frac{2-\mu}{2+\mu} + \sum_{j=1}^{\infty} \frac{(2-\mu)\mu^j}{(2+\mu)2^{j-1}} s^j.$$

Again, the reasonings above work only for $\mu < 2$, and for $\mu \geq 2$ we obtain that $\hat{r}(\mu) = 0$ and the quasispecies distribution is degenerate (see also Fig. 2).

Example 3.3. Let

$$\mathbf{r}_\infty = (1, 2, 0, \dots, 0, \dots).$$

Then $\mathbf{r}_\infty \circ P_\infty(s) = P_\infty(0) + 2P'_\infty(0)s$, and (2.5) takes the form

$$-\mu(1-s)P_\infty(s) + P_\infty(0) + 2P'_\infty(0)s = \hat{r}_\infty P_\infty(s).$$

We must assume that $P_\infty(0) = \hat{p}_{\infty,0} = 0$ since $\max r_k = 2 = r_1 > r_0$, therefore we cannot just plug $s = 0$ in the last equality. Instead, we differentiate it and find

$$-\mu(1-s)P'_\infty(s) + \mu P_\infty(s) + 2P'_\infty(0) = \hat{r}_\infty P'_\infty(s).$$

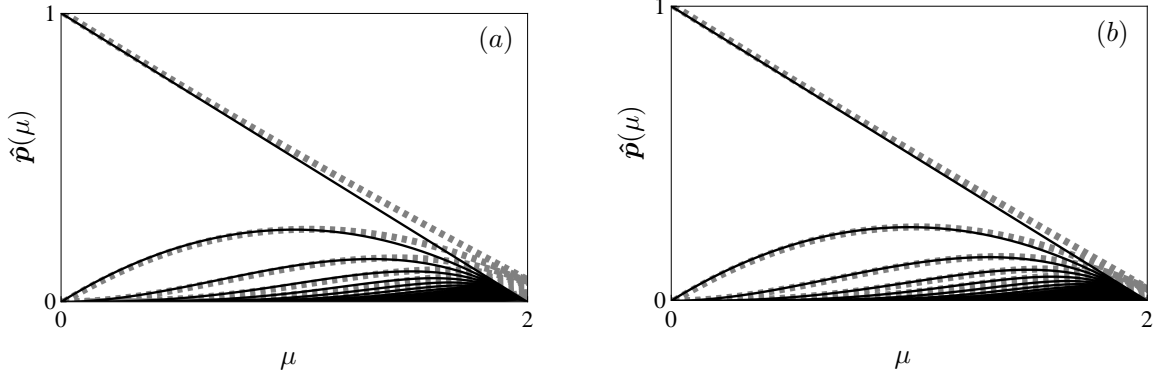


Figure 3: Comparison of numerical calculations for the fitness landscape in Example 3.3 with the theoretical predictions. The black solid lines are the exact solutions for the case $N \rightarrow \infty$ and the grey dashed lines are numerical computations. (a) $N = 50$; (b) $N = 100$

Then, for $s = 0$, assuming that $\hat{p}_{\infty,0} = 0$ and $P'_{\infty}(0) \neq 0$,

$$\hat{r}_{\infty} = 2 - \mu.$$

Therefore,

$$-\mu(1-s)P_{\infty}(s) + 2P'_{\infty}(0)s = (2-\mu)P_{\infty}(s),$$

or

$$P_{\infty}(s) = \frac{2P'_{\infty}(0)s}{2-\mu s},$$

which, together with $P_{\infty}(1) = 1$, gives

$$P_{\infty}(s) = \frac{(2-\mu)s}{2-\mu s} = \sum_{j=1}^{\infty} \left(1 - \frac{\mu}{2}\right) \frac{\mu^{j-1}}{2^{j-1}} s^j,$$

which holds only for $\mu < 2$, for $\mu \geq 2$ the distribution becomes degenerate (see Fig. 3).

Remark 3.4. In general, if $\max r_k = r_a$ and $r_j < r_a$ for $j > a$ then we need additional initial conditions

$$P_{\infty}(0) = P'_{\infty}(0) = \dots = P_{\infty}^{(a-1)}(0) = 0.$$

These initial conditions are motivated by the comparison of the theoretical predictions with the numerical calculations, and at this point we lack an analytical proof of the validity of these conditions in general.

Example 3.5 (A geometric landscape). Let $0 < q < 1$ and

$$\mathbf{r}_{\infty} = (1, q, q^2, \dots, q^N, \dots),$$

then $\mathbf{r}_{\infty} \circ P_{\infty}(s) = P_{\infty}(qs)$, and (2.5) reads

$$-\mu(1-s)P_{\infty}(s) + P_{\infty}(qs) = \hat{r}_{\infty}P_{\infty}(s).$$

Plugging $s = 0$ and assuming $P_\infty(0) \neq 0$ we find

$$\hat{r}_\infty = 1 - \mu,$$

which implies

$$P_\infty(qs) = (1 - \mu s)P_\infty(s). \quad (3.2)$$

Using the fact $P_\infty(1) = 1$ and plugging into the last expression $s = 1$, $s = q$, $s = q^2, \dots$ we find

$$P_\infty(q) = 1 - \mu, P_\infty(q^2) = (1 - \mu q)P_\infty(q) = (1 - \mu)(1 - \mu q), \dots, P_\infty(q^n) = \prod_{j=0}^{n-1} (1 - \mu q^j), \dots$$

Taking the limit $n \rightarrow \infty$ implies

$$P_\infty(0) = \lim_{n \rightarrow \infty} P_\infty(q^n) = \prod_{j=0}^{\infty} (1 - \mu q^j).$$

Assuming $P_\infty(s) = \sum_{j=0}^{\infty} \hat{p}_j s^j$ in (3.2), we obtain

$$q^n \hat{p}_n = \hat{p}_n - \mu \hat{p}_{n-1}, \quad \text{or} \quad \hat{p}_n = \frac{\mu}{1 - q^n} \hat{p}_{n-1},$$

which gives, for $0 < \mu < 1$, the limit distribution (see also Fig. 4)

$$\hat{p}_0 = \prod_{j=0}^{\infty} (1 - \mu q^j), \quad \hat{p}_n = \frac{\mu^n}{\prod_{k=1}^n (1 - q^k)} \hat{p}_0.$$

The value $\mu = 1$ is critical and corresponds to the error threshold. See Fig. 4 for comparison of the theoretical predictions and numerical computations.

We note that there are effective methods to calculate the expressions of the form $\prod_{j=0}^{\infty} (1 - \mu q^j)$ numerically. For instance, in *Mathematica*® this is done with the help of function `QPochhammer` $[\mu, q]$.

Example 3.6 (Additive or Fujiyama fitness landscape). The only fitness landscape for which problem (1.2) with the mutation scheme (1.3) can be analytically solved for the selection–mutation equilibrium is the additive fitness landscape, which we define here as

$$\mathbf{r} = \mathbf{r}_N = \left(1, 1 - \frac{1}{N}, 1 - \frac{2}{N}, \dots, 1 - \frac{N}{N} = 0\right).$$

The solution can be found in [3, 14] and arguably is most naturally derived using the tensor products for the representation of matrices $\mathbf{M} + \mathbf{J}\mathbf{M}$. Here we re-derive the same solution using the generating function approach.

Since

$$\mathbf{r}_N \circ P_N(s) = P_N(s) - \frac{s}{N} P'_N(s),$$

then (2.4) reads

$$\frac{\mu}{N} (1 - s^2) P'_N(s) - \mu (1 - s) P_N(s) + P_N(s) - \frac{s}{N} P'_N(s) = \hat{r}_N P_N(s).$$

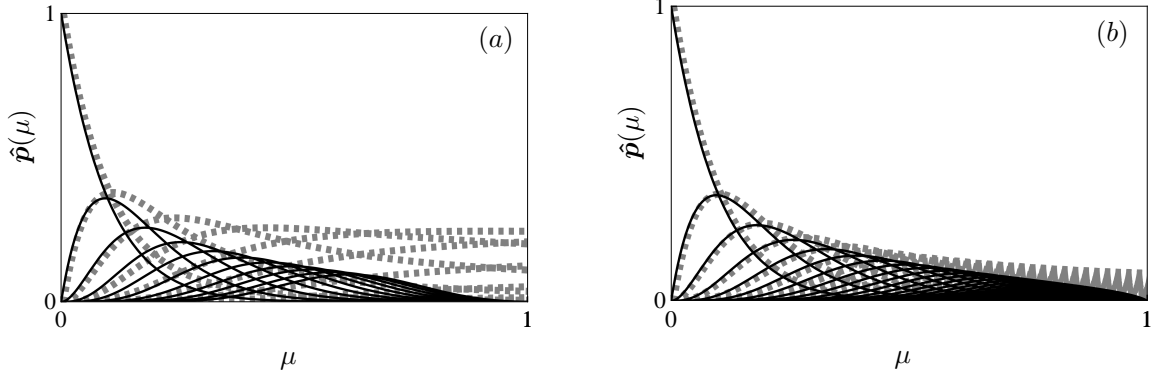


Figure 4: Comparison of numerical calculations for the fitness landscape in Example 3.5 with the theoretical predictions. The black solid lines are the exact solutions for the case $N \rightarrow \infty$ and the grey dashed lines are numerical computations. (a) $N = 10$; (b) $N = 50$

Making the substitution $P_N(s) = W^N(s)$ yields the ODE

$$\mu(1 - s^2)W'(s) - \mu(1 - s)W(s) + W(s) - sW'(s) = \hat{r}_N W(s),$$

or

$$\frac{W'(s)}{W(s)} = \frac{A}{s + a} + \frac{B}{s + b},$$

where

$$a = \frac{\sqrt{1 + 4\mu^2} + 1}{2\mu}, \quad b = \frac{\sqrt{1 + 4\mu^2} - 1}{2\mu}, \quad ab = 1, \quad a > 0, \quad 0 < b < 1,$$

$$A = \frac{1}{2} + \frac{2\hat{r}_N + 2\mu - 1}{2\sqrt{1 + 4\mu^2}}, \quad B = \frac{1}{2} - \frac{2\hat{r}_N + 2\mu - 1}{2\sqrt{1 + 4\mu^2}}, \quad A + B = 1.$$

Integrating this ODE yields

$$W(s) = C(s + a)^A(s + b)^B,$$

with the condition $C(1 + a)^A(1 + b)^B = 1$. We formally have that $W(b) = 0$, but this cannot occur, since $0 < b < 1$ and the polynomial $P_N(s)$ has all non-negative coefficients and is not equal to zero anywhere on the interval $[0, 1]$. This implies that $B = 0$, and therefore $A = 1$, and

$$W(s) = \frac{s + a}{1 + a}.$$

Moreover, the condition $B = 0$ implies that independently of N

$$\hat{r}_N = \frac{1 - 2\mu + \sqrt{1 + 4\mu^2}}{2}.$$

The final solution is

$$P_N(s) = \left(\frac{s + a}{1 + a} \right)^N,$$

and therefore the steady state distribution is binomial:

$$\hat{p}_{N,j} = \binom{N}{j} \frac{b^j}{(1+b)^N} \quad j = 0, \dots, N,$$

which holds for any $\mu > 0$, there exists no error threshold for this fitness landscape.

Example 3.7. Consider a close relative of the additive fitness landscape in the form

$$\mathbf{r}_\infty = \left(1, 1 - \frac{1}{K}, 1 - \frac{2}{K}, \dots, 1 - \frac{K}{K} = 0, \dots\right),$$

where $K \in \mathbf{N}$ and does not depend on N .

We have

$$\mathbf{r}_\infty \circ P_\infty(s) = \sum_{a=0}^{K-1} \left(1 - \frac{a}{K}\right) \hat{p}_a s^a,$$

and (2.5) reads

$$-\mu(1-s)P_\infty(s) + \sum_{a=0}^{K-1} \left(1 - \frac{a}{K}\right) \hat{p}_a s^a = \hat{\tau}_\infty P_\infty(s).$$

Plugging $s = 0$ and assuming $P_\infty(0) \neq 0$ implies

$$\hat{\tau}_\infty = 1 - \mu.$$

Therefore,

$$P_\infty(s) = \frac{\sum_{a=0}^{K-1} \left(1 - \frac{a}{K}\right) \hat{p}_a s^a}{1 - \mu s}, \quad (3.3)$$

and we need to determine \hat{p}_a for $a = 0, \dots, K-1$. Since we have

$$\sum_{a=0}^{K-1} \left(1 - \frac{a}{K}\right) \hat{p}_a s^a = (1 - \mu s) \sum_{a=0}^{\infty} \hat{p}_a s^a,$$

then for $1 \leq a \leq K-1$

$$\left(1 - \frac{a}{K}\right) \hat{p}_a = \hat{p}_a - \mu \hat{p}_{a-1},$$

or

$$\hat{p}_a = \frac{\mu K}{a} \hat{p}_{a-1}.$$

The last recurrent formula implies that for $1 \leq a \leq K-1$

$$\hat{p}_a = \frac{(\mu K)^a}{a!} \hat{p}_0.$$

To determine \hat{p}_0 , we use $P_\infty(1) = 1$, which yields, for $\mu < 1$,

$$\sum_{a=0}^{K-1} \left(1 - \frac{a}{K}\right) \hat{p}_a = 1 - \mu,$$

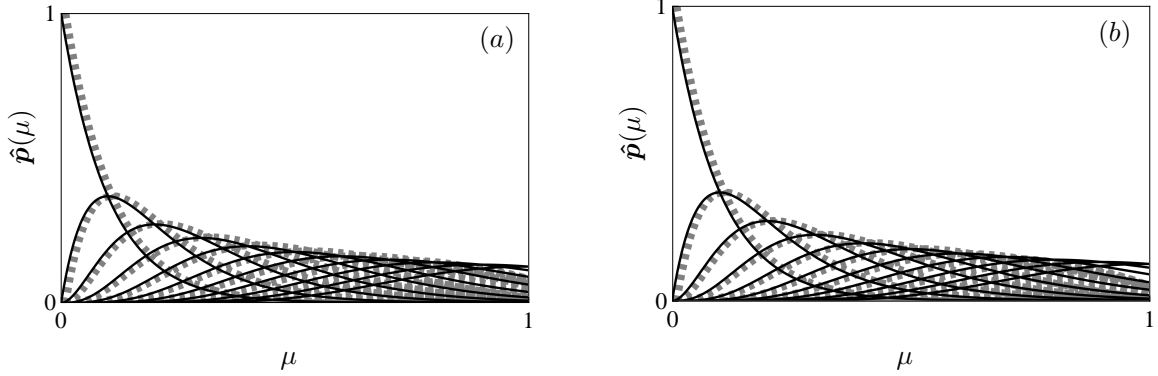


Figure 5: Comparison of numerical calculations for the fitness landscape in Example 3.7 with the theoretical predictions. The black solid lines are the exact solutions for the case $N \rightarrow \infty$ and the grey dashed lines are numerical computations. (a) $N = 100$, $K = 10$; (b) $N = 200$, $K = 10$

or, using the expressions for \hat{p}_a :

$$1 = \hat{p}_0 \left(\sum_{a=0}^{K-2} \frac{(\mu K)^a}{a!} + \frac{(\mu K)^{K-1}}{(1-\mu)(K-1)!} \right),$$

which allows us to find \hat{p}_0 . Now we determined all \hat{p}_a for $0 \leq a \leq K-1$ and from (3.3) we have that for $j \geq K$

$$\hat{p}_j = \sum_{a=0}^{K-1} \left(1 - \frac{a}{K} \right) \hat{p}_a \mu^{j-a}.$$

To simplify the expressions for \hat{p}_j we note that due to the central limit theorem (CLT), for $K \rightarrow \infty$,

$$\hat{p}_0 \sim e^{-\mu K}.$$

Indeed, if ξ_1, \dots, ξ_K are independent identically distributed Poisson random variables with the mean $\mathbb{E}\xi_i = \mu$, then $X_K = \xi_1 + \dots + \xi_K$ has the Poisson distribution with the mean and the variance $\mathbb{E}X_K = \text{Var } X_K = \mu K$. Therefore, using the CLT,

$$\begin{aligned} \sum_{a=0}^{K-2} \frac{(\mu K)^a}{a!} e^{-\mu K} + \frac{(\mu K)^{K-1}}{(1-\mu)(K-1)!} e^{-\mu K} &= \mathbb{P}(X_K \leq K-2) + \frac{\mathbb{P}(X_K = K-1)}{1-\mu} \\ &= \mathbb{P}\left(\frac{X_K - \mu K}{\sqrt{\mu K}} \leq \frac{(1-\mu)K-2}{\sqrt{\mu K}}\right) + \frac{\mathbb{P}(X_K = K-1)}{1-\mu} \rightarrow \mathbb{P}(X < \infty) + 0 = 1 \end{aligned}$$

as $K \rightarrow \infty$. Here X is the standard normally distributed random variable.

Therefore, for the case $1 \ll K \ll N$ the equilibrium distribution is approximately Poisson with parameter μK (see Fig. 5).

In all the examples above the limit mean fitness \hat{r}_∞ can be determined from the maximum principle (1.4). The next example shows that in the case when $r(x)$ has a point of discontinuity such that at this point function $r(x)$ neither left nor right continuous then a formal application of the maximum principle (1.4) may lead to incorrect conclusions.

Example 3.8. Let $N = 2A$ be an even number, and

$$\mathbf{r}_N = (0, \dots, 0, 1, 0, \dots, 0),$$

where 1 is exactly at the A -th position. In [6], using a parametric solution to the eigenvalue problem (this solution is outlined in Appendix A.1), it was proved that

$$\hat{r}_\infty = \sqrt{\mu^2 + 1} - \mu,$$

which is defined for any $\mu > 0$, there is no error threshold for this fitness landscape. In Appendix A.2 we re-derive this result using a new approach.

The maximum principle (1.4) for this example cannot be applied because $r(x)$ is neither left nor right continuous at the point $x = 0.5$. Formal application of the maximum principle leads to incorrect conclusion (e.g., for $\mu = 1$ it predicts that $\hat{r} \approx 1$, which is wrong, the exact value is $\sqrt{2} - 1$). In [6] also the expressions for the limit distribution $\hat{\mathbf{p}}_\infty$ were given without a full derivation. Here we show, using the method of generating functions, that the coordinates of the selection–mutation equilibrium indeed can be found in an explicit form.

In the following it will be convenient to consider the generating functions in the form of the Laurent series

$$U(s) = \sum_{n=-\infty}^{\infty} u_n s^n.$$

In (2.4) we make the substitution $P_N(s) = P_{2A}(s) = s^A U_A(s)$. Since the fitness landscape is symmetric, then the coefficients of $U_A(s)$ are also symmetric:

$$U_A(s) = u_{A,0} + \sum_{n=1}^A u_{A,n}(s^n + s^{-n}), \quad \hat{p}_{A \pm n} = u_{A,n}, \quad U_A(1) = 1.$$

After dividing by $2A$ equation (2.4) becomes

$$\frac{\mu}{2A}(1 - s^2)U'_A(s) + \frac{\mu}{2}(s^{-1} + s - 2)U_A(s) + u_{A,0} = \hat{r}_{2A}U_A(s).$$

Similarly to $\mathcal{H}1$ – $\mathcal{H}3$ we assume that, given that $A \rightarrow \infty$, the first term in the last equality vanishes, and $U_A(s)$ turns into

$$U_\infty(s) = u_0 + \sum_{n=1}^{\infty} u_n(s^n + s^{-n}), \quad U_\infty(1) = 1 = u_0 + 2 \sum_{n=1}^{\infty} u_n.$$

Therefore, we have the limit equation

$$\frac{\mu}{2}(s^{-1} + s - 2)U_\infty(s) + u_0 = \hat{r}_\infty U_\infty(s), \quad U_\infty(1) = 1,$$

or

$$u_0 = \left(\hat{r}_\infty + \mu - \frac{\mu}{2}(s^{-1} + s) \right) \left(u_0 + \sum_{n=1}^{\infty} (s^n + s^{-n}) \right).$$

Plugging in $s = 1$ we find

$$u_0 = \hat{r}_\infty.$$

Moreover, by equating the coefficients at $s^{-n} + s^n$, we obtain the system

$$u_0 = (\hat{r}_\infty + \mu)u_0 - \mu u_1, \quad 0 = (\hat{r}_\infty + \mu)u_n - \frac{\mu}{2}(u_{n-1} + u_{n+1}), \quad n \geq 1.$$

This system has the following solution, which can be directly checked,

$$u_n = \hat{r}_\infty \left(\frac{1 - \hat{r}_\infty}{1 + \hat{r}_\infty} \right)^n, \quad \hat{r}_\infty = \sqrt{\mu^2 + 1} - \mu, \quad n = 0, 1, \dots$$

Therefore, the limit distribution $\hat{\mathbf{p}}_\infty$ is two-sided geometric, and for large $N = 2A$ we have approximately

$$\hat{p}_{A \pm n} \approx u_n = \hat{r}_\infty \left(\frac{1 - \hat{r}_\infty}{1 + \hat{r}_\infty} \right)^n.$$

Numerical computations confirm this conclusion (see Fig. 6).

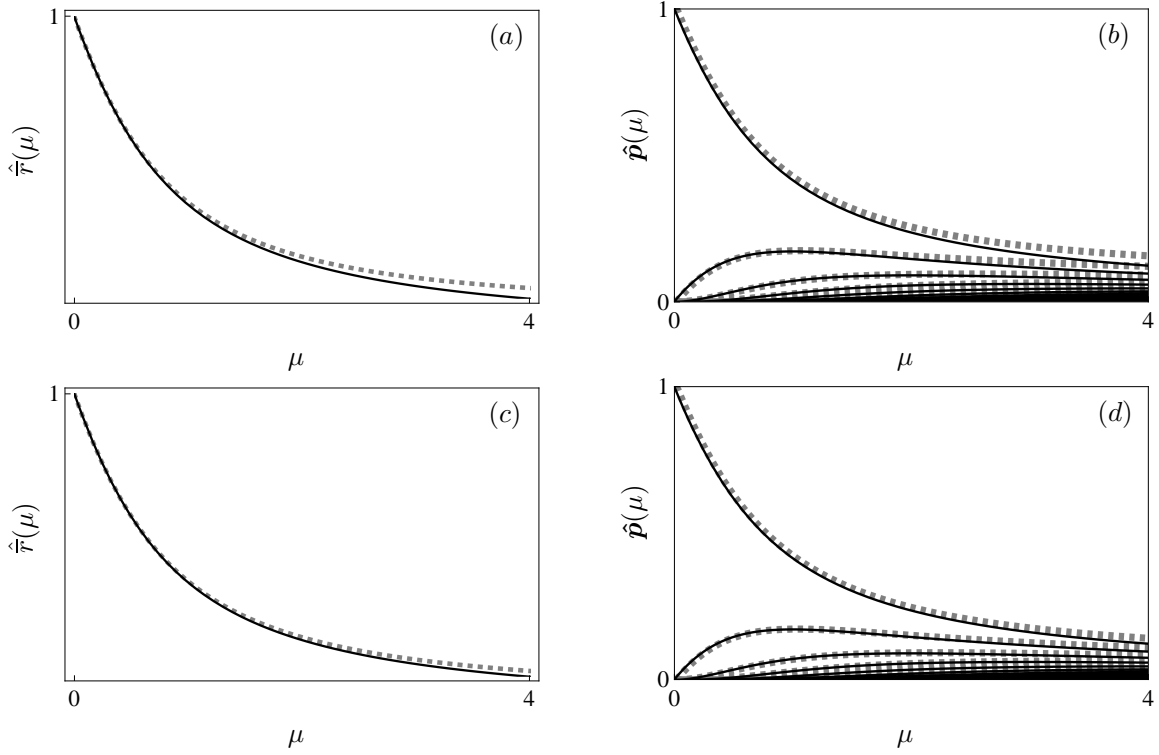


Figure 6: Comparison of numerical calculations for the fitness landscape in Example 3.8 with the theoretical predictions. The black solid lines are the exact solutions for the case $N \rightarrow \infty$ and the grey dashed lines are numerical computations. (a) and (b) show the mean population fitness and the quasispecies distribution respectively for $N = 100$; (c) and (d) show the same for $N = 200$

Example 3.9 (General formulas). As a final example, consider now a general fitness landscape, given by

$$\mathbf{r}_\infty = (r_0, r_1, \dots), \quad r_0 > r_i, \quad i = 1, 2, \dots \quad (3.4)$$

We can prove

Lemma 3.10. *Suppose (dropping the subscript ∞ for notational convenience) that*

$$P(s) = \sum_{n=0}^{\infty} \hat{p}_n s^n$$

gives a non-degenerate limit distribution

$$\hat{\mathbf{p}}_{\infty} = (\hat{p}_0, \hat{p}_1, \dots)$$

such that $\hat{p}_0 = P(0) > 0$.

Then

$$\hat{\tau}_{\infty} = r_0 - \mu, \quad \hat{p}_n = \frac{\mu^n \hat{p}_0}{\prod_{j=1}^n (r_0 - r_j)}, \quad n > 0, \quad \hat{p}_0 = \frac{1}{1 + \sum_{n=1}^{\infty} \frac{\mu^n}{\prod_{j=1}^n (r_0 - r_j)}}. \quad (3.5)$$

Expressions (3.5) generalize considered above Examples 3.1–3.3 and 3.5.

Proof. We rewrite (2.5) as follows:

$$\mu P(s) = \frac{\hat{\tau}_{\infty} P(s) - \mathbf{r}_{\infty} \circ P(s)}{s - 1} = \frac{\hat{\tau}_{\infty} (P(s) - 1) - (\mathbf{r}_{\infty} \circ P(s) - \hat{\tau}_{\infty})}{s - 1},$$

or, taking into account (2.6),

$$\mu \sum_{n=0}^{\infty} \hat{p}_n s^n = \hat{\tau}_{\infty} \sum_{n=1}^{\infty} \hat{p}_n \frac{s^n - 1}{s - 1} - \sum_{n=1}^{\infty} r_n \hat{p}_n \frac{s^n - 1}{s - 1},$$

or

$$\mu \sum_{n=0}^{\infty} \hat{p}_n s^n = \hat{\tau}_{\infty} \sum_{n=1}^{\infty} \hat{p}_n (1 + s + \dots + s^{n-1}) - \sum_{n=1}^{\infty} r_n \hat{p}_n (1 + s + \dots + s^{n-1}). \quad (3.6)$$

Substituting $s = 0$ yields

$$\mu \hat{p}_0 = \hat{\tau}_{\infty} \sum_{n=1}^{\infty} \hat{p}_n - \sum_{n=1}^{\infty} r_n \hat{p}_n = \hat{\tau}_{\infty} (1 - \hat{p}_0) - (\hat{\tau}_{\infty} - r_0 \hat{p}_0) = (r_0 - \hat{\tau}_{\infty}) \hat{p}_0.$$

By assumption $\hat{p}_0 > 0$, and we get the first equality in (3.5).

Now consider the coefficient at s in (3.6). We have

$$\mu \hat{p}_1 = \hat{\tau}_{\infty} (1 - \hat{p}_0 - \hat{p}_1) - (\hat{\tau}_{\infty} - r_0 \hat{p}_0 - r_1 \hat{p}_1) = (r_0 - \hat{\tau}_{\infty}) \hat{p}_0 - (\hat{\tau}_{\infty} - r_1) \hat{p}_1,$$

or, using the first inequality in (3.5),

$$\hat{p}_1 = \frac{\mu \hat{p}_0}{r_0 - r_1}.$$

Proceeding by induction on n and comparing the coefficients at s^n in (3.6), we prove the second equality in (3.5). The last equality follows from the condition $\sum_{n=0}^{\infty} \hat{p}_n = 1$. \blacksquare

Remark 3.11. The expressions for the limit distribution in Lemma 3.10 can be generalized to the case

$$\mathbf{r}_\infty = (r_0, r_1, \dots), \quad r_k \geq r_i, \quad i = 0, 1, \dots, k-1, \quad r_k > r_i, \quad i = k+1, k+2, \dots$$

If one assumes additionally (as the numerical evidence suggests, see also Remark 3.4) that

$$\hat{p}_0 = \dots = \hat{p}_{k-1} = 0, \quad \hat{p}_k > 0,$$

then

$$\hat{r}_\infty = r_k - \mu, \quad \hat{p}_{k+n} = \frac{\mu^n \hat{p}_k}{\prod_{j=1}^n (r_k - r_{k+j})}, \quad n > 0, \quad \hat{p}_k = \frac{1}{1 + \sum_{n=1}^{\infty} \frac{\mu^n}{\prod_{j=1}^n (r_k - r_{k+j})}}. \quad (3.7)$$

Corollary 3.12. Assume that the conditions of Lemma 3.10 hold. If, additionally, the limit fitness landscape is such that $r_n = 0$ for $n > n_0$ then the generating function $P(\infty)$ is rational and the limit distribution is asymptotically geometric.

Corollary 3.13. Assume that the conditions of Lemma 3.10 hold. Then formulas (3.7) provide a solution for the non-degenerate limit distribution if and only if

$$\mu < \mu^* = \liminf_n \sqrt[n]{\prod_{j=1}^n (r_k - r_{k+j})} \leq r_k. \quad (3.8)$$

Proof. Indeed, from (3.7), the necessary and sufficient condition for the non-degenerate distribution to exist is the convergence of the series

$$\sum_{n=1}^{\infty} \frac{\mu^n}{\prod_{j=1}^n (r_k - r_{k+j})},$$

which, by the Cauchy–Hadamard formula, converges if (3.8) holds and diverges if $\mu > \mu^*$. ■

Remark 3.14. The critical value μ^* of the mutation rate in (3.8) should be considered the threshold value of the error threshold.

To illustrate how these general expressions work, consider the fitness landscape

$$\mathbf{r}_\infty = (2, 0, 1, 0, 1, 0, 1, 0, \dots). \quad (3.9)$$

Equation (3.8) predicts that for $\mu < \mu^* = \sqrt{2}$ we must have

$$\hat{r}_\infty = 2 - \mu,$$

and, for several first coordinates of the equilibrium distribution

$$\hat{p}_0 = \frac{2 - \mu^2}{2 + \mu}, \quad \hat{p}_1 = \frac{2 - \mu^2}{2 + \mu} \cdot \frac{\mu}{2}, \quad \hat{p}_2 = \frac{2 - \mu^2}{2 + \mu} \cdot \frac{\mu^2}{2}, \quad \hat{p}_3 = \frac{2 - \mu^2}{2 + \mu} \cdot \frac{\mu^3}{4}.$$

Comparison of these exact theoretical predictions with numerical computations are given in Fig. 7.

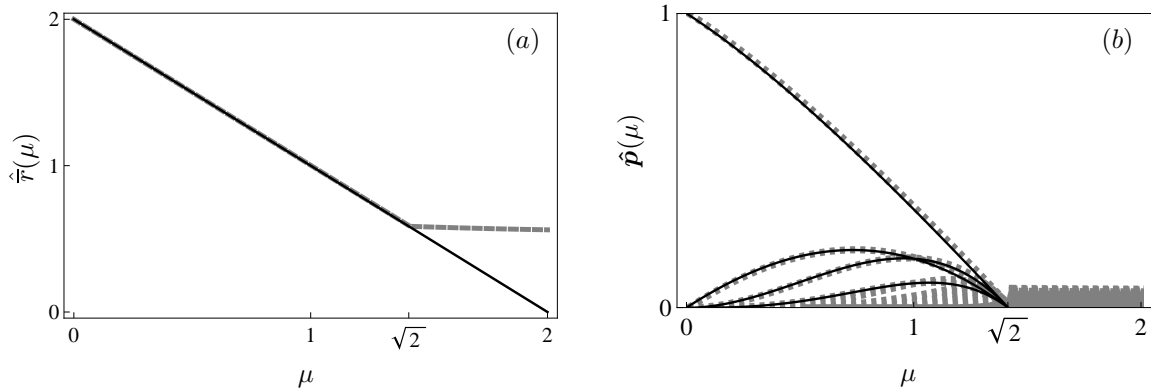


Figure 7: Comparison of numerical calculations for the fitness landscape (3.9) in Example 3.9 with the theoretical predictions. The black solid lines are the exact solutions for the case $N \rightarrow \infty$ and the grey dashed lines are numerical computations for $N = 200$. (a) Mean population fitness. (b) The quasispecies distribution

4 Concluding remarks

We presented an analytical approach to calculate the mutation–selection equilibrium in the Crow–Kimura evolutionary model, which is based on the reformulation of the original eigenvalue problem as a nonlinear functional–differential equation for the unknown probability generating function and on taking a formal limit $N \rightarrow \infty$ for the sequence length. This approach provides closed analytical solutions for at least several special fitness landscapes, as we amply illustrated in the previous section. We remark that, to the best of our knowledge, in the existing literature only for two fitness landscapes these equilibrium distributions were written down explicitly, and most attention was concentrated on finding analytical expressions for the mean population fitness (the leading eigenvalue) and for some other population averages. With the advent of sequencing technique, as everyone witnessed for the last two decades, it is now completely feasible to sequence the whole population of quasispecies, and therefore our formulas can be used to further relate theory and experiment in the evolutionary questions.

While the approach suggested in Section 2 clearly works for all the examples we considered, the conditions are quite difficult to rigorously check and their proof constitutes an independent and deep problem. Our experience tells us that it is quite unlikely to present general (necessary and/or sufficient) conditions, which are easy to check, that would guarantee that our formal limit yields the correct result for an arbitrary fitness landscape. As of now each special case has to be tackled on its own, as we demonstrated for the single peaked landscape (Example 3.1) in [6]. Probably a first realistic step in the search of the general conditions is to prove that the formulas in Example 3.9 follow rigorously from the assumption on the unique fixed maximum of the fitness landscape. We conclude our text with this open problem.

A Parametric solutions to the basic eigenvalue problem

The suggested general approach of the generating functions is heuristic. Rigorous proofs can be obtained using the parametric solution method introduced in [6]. In this appendix we give a concise form for the method used in [6] (see Appendix A.1) and also introduce a new parametric solution in Appendix A.2, which is used to prove the result from Example 3.8 on the limit form of \hat{r}_∞ .

Recall that we are interested in finding the dominant eigenvalue and the corresponding positive eigenvector of the problem

$$(\mathbf{M} + \mu \mathbf{Q})\hat{\mathbf{p}} = \hat{m} \hat{\mathbf{p}},$$

where $\mu \geq 0$, $\mathbf{M} = \text{diag}(m_0, \dots, m_N)$, and the matrix \mathbf{Q} has the form

$$\mathbf{Q} = \begin{bmatrix} -N & 1 & 0 & 0 & \dots & \dots & 0 \\ N & -N & 2 & 0 & \dots & \dots & 0 \\ 0 & N-1 & -N & 3 & \dots & \dots & 0 \\ 0 & 0 & N-2 & -N & \dots & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & 0 \\ 0 & 0 & \dots & \dots & 2 & -N & N \\ 0 & 0 & \dots & \dots & 0 & 1 & -N \end{bmatrix}.$$

We introduce the notations

$$\mathbf{S} = \frac{1}{N} \mathbf{Q}, \quad \mathbf{R} = \frac{1}{N} \mathbf{M}, \quad \hat{r} = \frac{1}{N} \hat{m}.$$

Then the original eigenvalue problem takes the form

$$(\mathbf{R} + \mu \mathbf{S})\hat{\mathbf{p}} = \hat{r} \hat{\mathbf{p}},$$

or, after introducing a new parameter $u = \mu/\hat{r}$,

$$\frac{1}{\hat{r}} \mathbf{R}\hat{\mathbf{p}} + u \mathbf{S}\hat{\mathbf{p}} = \hat{\mathbf{p}}. \tag{A.1}$$

A.1 Approach A

In [6] it was shown that

$$\mathbf{C}^{-1} \mathbf{Q} \mathbf{C} = -2 \text{diag}(0, 1, 2, \dots, N),$$

where $\mathbf{C} = (c_{ka})$ is the matrix composed (by columns) of the coefficients of the generating polynomials

$$P_a(s) = \sum_{k=0}^N c_{ka} s^k = (1-s)^a (1+s)^{N-a},$$

and possesses the property $\mathbf{C}^2 = 2^N \mathbf{I}$, where \mathbf{I} is the identity matrix. Using this information, we have

$$\mathbf{C}^{-1} \mathbf{S} \mathbf{C} = -2 \text{diag}\left(0, \frac{1}{N}, \frac{2}{N}, \dots, 1\right).$$

Equation (A.1) can be written as

$$\hat{r} \hat{\mathbf{p}} = (\mathbf{I} - u\mathbf{S})^{-1} \mathbf{R} \hat{\mathbf{p}}.$$

Using

$$(\mathbf{I} - u\mathbf{S})^{-1} = \mathbf{C}^{-1}(\mathbf{I} - u\mathbf{C}^{-1}\mathbf{S}\mathbf{C})^{-1}\mathbf{C} = \mathbf{C}^{-1} \text{diag} \left(1, \frac{1}{1 + \frac{2u}{N}}, \frac{1}{1 + \frac{4u}{N}}, \dots, \frac{1}{1 + \frac{2uN}{N}} \right) \mathbf{C},$$

we find

$$(\mathbf{I} - u\mathbf{S})^{-1} = \mathbf{F}(u) = (F_{ab}(u)), \quad F_{ab} = \frac{1}{2^N} \sum_{k=0}^N \frac{c_{ak}c_{kb}}{1 + \frac{2ku}{N}}.$$

Therefore, \hat{r} is the dominant eigenvalue of the matrix $\mathbf{F}(u)\mathbf{R}$, and $\hat{\mathbf{p}}$ is the corresponding eigenvector, both of which can be represented in the parametric form, depending on parameter u . The details how to write down the explicit formulas depending on the number of nonzero elements of the vector $\mathbf{m} = (m_0, \dots, m_N)$ are given in [6].

A.2 Approach B

The equality

$$\hat{r} \hat{\mathbf{p}} = (\mathbf{I} - u\mathbf{S})^{-1} \mathbf{R} \hat{\mathbf{p}}$$

can be written as

$$\hat{r} \hat{\mathbf{p}} = ((1 + u)\mathbf{I} - u(\mathbf{S} + \mathbf{I}))^{-1} \mathbf{R} \hat{\mathbf{p}}. \quad (\text{A.2})$$

Let

$$\mathbf{B} = \mathbf{S} + \mathbf{I} = \begin{bmatrix} 0 & 1/N & 0 & 0 & \dots & \dots & 0 \\ 1 & 0 & 2/N & 0 & \dots & \dots & 0 \\ 0 & 1 - 1/N & 0 & 3/N & \dots & \dots & 0 \\ 0 & 0 & 1 - 2/N & 0 & \dots & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \dots & 2/N & 0 & 1 \\ 0 & 0 & \dots & \dots & 0 & 1/N & 0 \end{bmatrix},$$

i.e., the matrix \mathbf{B} is two diagonal, stochastic, and such that

$$\mathbf{C}^{-1}\mathbf{B}\mathbf{C} = \text{diag} \left(1, 1 - \frac{2}{N}, 1 - \frac{4}{N}, \dots, -1 \right).$$

Direct observations yield that all natural powers of \mathbf{B} also will be stochastic. For the even powers $\mathbf{B}^{2k} = (b_{ij}^{(2k)})$ we have $b_{ij}^{(2k)} = 0$ if $i + j$ is odd, and for the odd powers $\mathbf{B}^{2k-1} = (b_{ij}^{(2k-1)})$ we have $b_{ij}^{(2k-1)} = 0$ if $i + j$ is even. Moreover, the properties of \mathbf{B} and \mathbf{C} imply

$$\begin{aligned} \lim_{k \rightarrow \infty} \mathbf{B}^{2k} &= \mathbf{C} \text{diag}(1, 0, \dots, 0, 1) \mathbf{C}^{-1} = \frac{1}{2^N} \left((1 + (-1)^{i+j}) \binom{N}{i} \right), \\ \lim_{k \rightarrow \infty} \mathbf{B}^{2k-1} &= \mathbf{C} \text{diag}(1, 0, \dots, 0, -1) \mathbf{C}^{-1} = \frac{1}{2^N} \left((1 - (-1)^{i+j}) \binom{N}{i} \right). \end{aligned}$$

Using the introduced notations equality (A.2) can be written as

$$\hat{\mathbf{r}} \hat{\mathbf{p}} = (1 + u)^{-1} \left(\mathbf{I} - \frac{u}{1 + u} \mathbf{B} \right)^{-1} \mathbf{R} \hat{\mathbf{p}} = \sum_{k=0}^{\infty} \frac{u^k}{(1 + u)^{k+1}} \mathbf{B}^k \mathbf{R} \hat{\mathbf{p}}, \quad (\text{A.3})$$

where the application of the geometric series can be justified by noting that the 1-norm of the matrix $\frac{u}{1+u} \mathbf{B}$ is less than 1 for any $u > 0$.

To illustrate how the parametric solution (A.3) works consider the case of the fitness landscape in Example 3.8. Let $N = 2A$ be even,

$$\mathbf{r} = (0, \dots, 0, 1, 0, \dots, 0),$$

where one is at the A -th position. In this case

$$\hat{\mathbf{r}} = \hat{p}_A,$$

and (A.3) takes the form

$$\hat{\mathbf{r}} \hat{p}_A = \sum_{k=0}^{\infty} \frac{u^k}{(1 + u)^{k+1}} b_{A,A}^{(k)} \hat{p}_A.$$

Since $b_{A,A}^{(k)} = 0$ for odd k and $\hat{\mathbf{r}} = \hat{p}_A > 0$ then

$$\hat{\mathbf{r}} = \sum_{k=0}^{\infty} \frac{u^{2k}}{(1 + u)^{2k+1}} b_{A,A}^{(2k)}.$$

From the properties of \mathbf{B} and \mathbf{C} we have

$$b_{A,A}^{(2k)} = \frac{1}{2^{2A}} \sum_{i=0}^{2A} c_{Ai} c_{iA} \left(1 - \frac{i}{A} \right)^{2k}.$$

Using the property that $c_{ij} \binom{N}{j} = c_{ji} \binom{N}{i}$ (see [6] for an easy proof) and the fact that c_{iA} are the coefficients of the polynomial

$$P_A(s) = \sum_{l=0}^A (-1)^l \binom{A}{l} s^{2l},$$

we find

$$c_{Ai} c_{iA} = \begin{cases} 0, & i = 2l + 1, \\ \binom{2l}{l} \binom{2(A-l)}{A-l}, & i = 2l, \end{cases}$$

and therefore

$$b_{A,A}^{(2k)} = \frac{1}{2^{2A}} \sum_{l=0}^A \binom{2l}{l} \binom{2(A-l)}{A-l} \left(1 - \frac{2l}{A} \right)^{2k}.$$

Using the approximation

$$\frac{1}{2^{2n}} \binom{2n}{n} \approx \frac{1}{\sqrt{\pi n}},$$

we obtain

$$b_{A,A}^{(2k)} \approx \frac{2}{2^{2A}} \binom{2A}{A} + \frac{1}{\pi} \sum_{l=1}^{A-1} \frac{\left(1 - \frac{2l}{A}\right)^{2k}}{\sqrt{l(A-l)}},$$

or, after taking the limit $A \rightarrow \infty$,

$$\lim_{A \rightarrow \infty} b_{A,A}^{(2k)} = \frac{1}{\pi} \int_0^1 \frac{(1-2x)^{2k} dx}{\sqrt{x(1-x)}} = \frac{2}{\pi} \int_0^{\pi/2} \cos^{2k} z dz = \frac{1}{2^{2k}} \binom{2k}{k}.$$

Therefore, in the limit of the infinite sequence length

$$\hat{r}_\infty = \sum_{k=0}^{\infty} \frac{u^{2k}}{(1+u)^{2k+1}} \frac{1}{2^{2k}} \binom{2k}{k} = \frac{1}{1+u} \frac{1}{\sqrt{1 - \frac{u^2}{(1+u)^2}}} = \frac{1}{\sqrt{2u+1}},$$

where we used the fact that

$$\frac{1}{\sqrt{1-x^2}} = \sum_{k=0}^{\infty} \frac{1}{2^{2k}} \binom{2k}{k} x^{2k}, \quad |x| < 1.$$

Finally, remembering that $u = \frac{\mu}{\hat{r}_\infty}$ gives

$$\hat{r}_\infty = \frac{1}{\sqrt{\frac{2\mu}{\hat{r}_\infty} + 1}},$$

from where

$$\hat{r}_\infty = \sqrt{1 + \mu^2} - \mu,$$

as it was stated in Example 3.8 and proved in [6] using, essentially, approach from Appendix A.1.

Acknowledgements: ASN's research is supported in part by ND EPSCoR and NSF grant #EPS-0814442.

References

- [1] E. Baake and W. Gabriel. Biological evolution through mutation, selection, and drift: An introductory review. In D. Stauffer, editor, *Annual Reviews of Computational Physics VII*, pages 203–264. World Scientific, 1999.
- [2] E. Baake and H.-O. Georgii. Mutation, selection, and ancestry in branching models: a variational approach. *Journal of Mathematical Biology*, 54(2):257–303, Feb 2007.
- [3] E. Baake and H. Wagner. Mutation–selection models solved exactly with methods of statistical mechanics. *Genetical research*, 78(1):93–117, 2001.
- [4] N. T. J. Bailey. *The elements of stochastic processes with applications to the natural sciences*, volume 25. John Wiley & Sons, 1990.

- [5] C. K. Biebricher and M. Eigen. The error threshold. *Virus research*, 107(2):117–127, 2005.
- [6] A. S. Bratus, A. S. Novozhilov, and Y. S. Semenov. Linear algebra of the permutation invariant Crow–Kimura model of prebiotic evolution. *Mathematical Biosciences*, 256:42–57, 2014.
- [7] R. Bürger. *The mathematical theory of selection, mutation, and recombination*. Wiley, 2000.
- [8] J. F. Crow and M. Kimura. *An introduction to population genetics theory*. New York, Evanston and London: Harper & Row, Publishers, 1970.
- [9] M. Eigen. Selforganization of matter and the evolution of biological macromolecules. *Naturwissenschaften*, 58(10):465–523, 1971.
- [10] M. Eigen, J. McCaskill, and P. Schuster. Molecular quasi-species. *Journal of Physical Chemistry*, 92(24):6881–6891, 1988.
- [11] M. Eigen and P. Schuster. The hypercycle. A principle of natural self-organization. Part A: Emergence of the hypercycle. *Naturwissenschaften*, 64(11):541–565, Nov 1977.
- [12] S. Galluccio. Exact solution of the quasispecies model in a sharply peaked fitness landscape. *Physical Review E*, 56(4):4526, 1997.
- [13] J. Hermisson, O. Redner, H. Wagner, and E. Baake. Mutation-selection balance: ancestry, load, and maximum principle. *Theoretical Population Biology*, 62(1):9–46, Aug 2002.
- [14] P. G. Higgs. Error thresholds and stationary mutant distributions in multi-locus diploid genetics models. *Genetical Research*, 63(01):63–78, 1994.
- [15] K. Jain and J. Krug. Adaptation in Simple and Complex Fitness Landscapes. In U. Bastolla, M. Porto, H. Eduardo Roman, and M. Vendruscolo, editors, *Structural approaches to sequence evolution*, chapter 14, pages 299–339. Springer, 2007.
- [16] D. S. Rumschitzki. Spectral properties of Eigen evolution matrices. *Journal of Mathematical Biology*, 24(6):667–680, 1987.
- [17] D. B. Saakian. A new method for the solution of models of biological evolution: Derivation of exact steady-state distributions. *Journal of Statistical Physics*, 128(3):781–798, 2007.
- [18] D. B. Saakian, C.-K. Hu, and H. Khachatryan. Solvable biological evolution models with general fitness functions and multiple mutations in parallel mutation-selection scheme. *Physical Review E*, 70(4):041908, 2004.
- [19] Y. S. Semenov, A. S. Bratus, and A. S. Novozhilov. On the behavior of the leading eigenvalue of the Eigen evolutionary matrices. *Mathematical Biosciences*, 258:134–147, 2014.
- [20] J. Swetina and P. Schuster. Self-replication with errors: A model for polynucleotide replication. *Biophysical Chemistry*, 16(4):329–345, 1982.

- [21] C. J. Thompson and J. L. McBride. On Eigen's theory of the self-organization of matter and the evolution of biological macromolecules. *Mathematical Biosciences*, 21(1):127–142, 1974.
- [22] G. P. Wagner and P. Krall. What is the difference between models of error thresholds and Muller's ratchet? *Journal of Mathematical Biology*, 32(1):33–44, 1993.
- [23] T. Wiehe. Model dependency of error thresholds: the role of fitness functions and contrasts between the finite and infinite sites models. *Genetical research*, 69(02):127–136, 1997.
- [24] C. O. Wilke. Quasispecies theory in the context of population genetics. *BMC Evolutionary Biology*, 5(1):44, 2005.